

1 Population genomics of *Wolbachia* and mtDNA in *Drosophila simulans* from California

2

3 Sarah Signor*

4

5

6 *Department of Molecular and Computational Biology, University of Southern California, Los Angeles,

7 California USA

8 Key words: *Drosophila simulans*, population genetics, *Wolbachia*

9

10 Communicating Author: Sarah Signor

11 Section of Molecular and Computational Biology

12 University of Southern California

13 Los Angeles, California USA

14 ssignor@usc.edu

15 Phone 213-740-3065

16

17

18

19

20

21

22 ABSTRACT: *Wolbachia pipientis* is an intracellular endosymbiont infecting many arthropods and filarial
23 nematodes. Little is known about the short-term evolution of *Wolbachia* or its interaction with its host.
24 *Wolbachia* is maternally inherited, resulting in co-inheritance of mitochondrial organelles such as
25 mtDNA. Here I explore the short-term evolution of *Wolbachia*, and the relationship between *Wolbachia*
26 and mtDNA, using a large inbred panel of *Drosophila simulans* infected with the *Wolbachia* strain wRi. I
27 find reduced diversity relative to expectation in both *Wolbachia* and mtDNA, but only mtDNA shows
28 evidence of a recent selective sweep or population bottleneck. I estimate *Wolbachia* and mtDNA titre in
29 each genotype, and I find considerable variation in both phenotypes, despite low genetic diversity in
30 *Wolbachia* and mtDNA. A phylogeny of *Wolbachia* and of mtDNA show that both trees are largely
31 unresolved, suggesting a recent origin of the infection derived from a single origin. Using *Wolbachia* and
32 mtDNA titre as a phenotype, I perform an association analysis with the nuclear genome and find several
33 regions implicated in the phenotype, including one which contains four CAAX-box protein processing
34 genes. CAAX-box protein processing can be an important part of host-pathogen interactions in other
35 systems, suggesting interesting directions for future research.

36

37

38

39

40

41

42

43 INTRODUCTION:

44 Heritable symbiotic associations such as that between *Drosophila* and *Wolbachia pipientis* have
45 widespread impact on host ecology and evolution. These types of heritable endosymbiotic relationships
46 are recognized as key drivers of evolution, but the intraspecific variation that effects their short-term
47 evolution is not well explored. *Wolbachia* are α -proteobacterial endosymbionts found in up to 40% of all
48 arthropod species¹⁻³. *Wolbachia* are maternally transmitted and spread through manipulating the
49 reproductive strategies of their host, using mechanisms such as feminization, male-killing, or cytoplasmic
50 incompatibility. The most common of these is cytoplasmic incompatibility, where mating between males
51 and females of the same species results in embryonic mortality if they have different *Wolbachia* infection
52 status⁴⁻⁸. *Wolbachia* may also confer certain protections upon their host, such as increased resistance to
53 certain viruses, or increased survival when exposed to certain environmental stressors⁹⁻¹⁵. *Wolbachia* is
54 one of the most abundant obligate intracellular parasites, given that 85% of animal species are insects.
55 This has profound meaning for evolutionary processes such as sexual selection and speciation^{16,17}.

56 *Wolbachia* strain wRi is known to have spread recently in the sister species to the model organism
57 *Drosophila melanogaster*, *D. simulans*^{4,5,8,18}. It was at ~95-100% frequency in Southern California
58 populations at the time its original sampling in the 1980's^{19,20}. It likely invaded California less than 25
59 years before it was first detected in 1984²¹. It is now thought to have been horizontally transmitted to *D.*
60 *simulans* from *D. ananassae*, though the same strain is also found in *D. suzukii*²². The maternal
61 transmission of *Wolbachia* means that as the microorganism spreads all maternally inherited organelles
62 spread along with it. Most notably mtDNA will be forced through a bottleneck, lowering the diversity of
63 mtDNA in infected populations^{18,23,24}. This will cause mtDNA and *Wolbachia* to be more closely
64 associated than nuclear genes, and this coupling has been demonstrated previously in *D. simulans*^{18,25,26}. In
65 fact, *D. simulans* is known to have three major mitochondrial haplotypes (*siI*, *siII*, and *siIII*) and two
66 subtypes (*siIIA* and *siIIB*) that harbor very little variation and that appear to be nonrandomly associated
67 with *Wolbachia* strains²⁷⁻²⁹. These mitochondrial haplotypes are largely allopatric, except for the presence
68 of both *siII* and *siIII* in Madagascar and La Reunion³⁰.

69 In *D. melanogaster* variation in *Wolbachia* has been well investigated using genomic data, though
70 this has not been done in the genomic era in *D. simulans*³¹. They found long lived associations between
71 mitochondrial and *Wolbachia* haplotypes and strong geographic structuring among cytotypes³¹⁻³³. This
72 study also observed that *Wolbachia* titre varied among fly populations as the result of intraspecific
73 nuclear genetic variation³¹. However, the assumption that it was due to intraspecific nuclear background
74 was based on the presence of a constant environment and no polymorphisms were identified that could be
75 affecting this phenotype. Very little is known about how *Wolbachia* interact with their hosts, though
76 recent work has uncovered evidence that deubiquitylating enzymes produced by *Wolbachia* and secreted
77 into the host cytoplasm mediate cytoplasmic incompatibility³⁴. *Wolbachia* DNA is also frequently
78 inserted into the host genome, though this has not occurred with *wRi* in *D. simulans*²¹. Genes involved in
79 the formation of germline stem cells such as *benign gonial cell neoplasm* and *bag-of-marbles* are
80 considered candidates for interacting with *Wolbachia*, and have been found to have unusual population
81 genetic patterns in *D. melanogaster*^{35,36}. *bag-of-marbles* has been suggested to interact with *Wolbachia*
82 due to fertility rescue in hypomorphs, but the interaction of this gene with *Wolbachia* in natural
83 populations is not clear^{22,35-37}. Notably, *Wolbachia* localizes in tissues differently depending upon the
84 strain and species so the interactions between the host and *Wolbachia* are likely to also be different^{38,39}.

85 *Wolbachia* infections must be maintained in host populations through transovarial transmission,
86 wherein *Wolbachia* is present in the germline at sufficient copy number to ensure transmission but not to
87 cause host pathology⁴⁰. *Wolbachia* titre has been shown to have important phenotypic effects on the
88 host^{11,41-49}. However, control of *Wolbachia* replication is not well understood, nor is the dependence of
89 this control on host background versus bacterial genotype^{11,50-52}. Differences in *Wolbachia* titre when it is
90 transinfected between species suggests a role of host background in controlling copy number, population
91 genomics in *D. melanogaster* suggest an effect of host background, and there does seem to be host-
92 specific patterns of tissue colonization⁵³⁻⁵⁵. However, multiple *Wolbachia* genotypes can also behave
93 differently in the same genetic background suggesting contributions from the bacterial genome^{51,56}. It is

94 also possible to select for greater *Wolbachia* densities, though the heritability of this is unclear^{57,58}.

95 Here I investigate the dynamics of *Wolbachia* and mtDNA in a large panel of *D. simulans* from a
96 single Californian population. I determine infection status of *Wolbachia* in the panel of *D. simulans*
97 genotypes. I look for signatures of selection in both genomes using summary statistics Tajima's D and π
98 and find that while *Wolbachia* patterns of variation are not unusual given its demographic history the
99 reduction in mtDNA diversity is suggestive of a recent bottleneck due either to selection or changes in
100 population size. I also measure linkage disequilibrium between mtDNA and *Wolbachia* as a proxy for co-
101 inheritance. Using whole genome sequences, I investigate the phylogeny of both *Wolbachia* and mtDNA
102 and find that in this population they are essentially unresolved. I investigate variation in the copy number
103 of both *Wolbachia* and mtDNA in this population using relative estimates derived from illumina
104 sequencing coverage compared to nuclear coverage. I find considerable copy number variation in this
105 population, and an association analysis using this as a phenotype implicates several genomic regions
106 potentially involved in mediating this phenotype. This includes a region containing multiple genes
107 involved CAAX-box protein prenylation, a process that is important for mediating the relationship
108 between host and pathogen in other systems⁵⁹⁻⁶¹.

109 METHODS

110 *Drosophila* strains

111 Strains are as described in⁶². Briefly, the 167 *D. simulans* lines were collected in the Zuma organic orchard in
112 Zuma beach, California in February of 2012 from a single pile of fermenting strawberries. Single mated females
113 were collected and inbred by 15 generations of full sib mating of their progeny. *Drosophila* were raised at a
114 constant temperature of 20° C with 12-hour light/dark cycles. They were raised on a standard glucose/yeast media,
115 and each library was constructed from adult females of similar age (less than one week).

116 **Data sources and processing**

117 The sequencing reads were downloaded from the NCBI Short Read Archive from project SRP075682.

118 Libraries were assembled using BWA mem (v. 0.7.5), and processed with samtools (v. 0.1.19) using
119 default parameters^{63,64}. The *Wolbachia* reference is the *w*Ri strain previously identified in Southern
120 California (Accession number NC_012416)²¹. The mtDNA reference is from *D. simulans* *w*⁵⁰¹, which is
121 haplogroup *siII* as expected for *D. simulans* from California (Accession number KC244284)⁶⁵. PCR
122 duplicates were removed using Picard MarkDuplicates (v. 2.9.4) and GATK (v. 3.7) was used for indel
123 realignment and SNP calling using default parameters (<http://picard.sourceforge.net>)⁶⁶. SNPs were called
124 jointly for all genotypes using Haplotypecaller⁶⁶. Individual consensus fasta sequences were produced
125 using SelectVariants to create individual vcf files and FastaAlternateReferenceMaker. Vcf files were
126 filtered for indels and non-biallelic SNPs using VCFtools (v. 0.1.13)⁶⁷. The files were also filtered for
127 SNPs with more than 10% missing data. The *Wolbachia* genome was filtered for regions of unusual
128 coverage or SNP density, for example two regions of the *Wolbachia* genome harbored ~40 SNPs within
129 two kb, far above background levels of variation (Supp. Fig. 1). These two regions coincided with regions
130 of unusually high coverage suggesting they are repeated elements.

131 **Prediction of *Wolbachia* infection status**

132 *Wolbachia* infection status was determined by calculating the mean depth of coverage of the assembly
133 and the breadth of coverage of the consensus sequence using bedtools⁶⁸. Depth of coverage refers to the
134 average read depth across the *Wolbachia* genome, while breadth of coverage refers to the number of bases
135 covered by at least two reads. Depth of coverage at each nucleotide was estimated using the genomecov
136 function, while breadth was estimated using the coverage function. Predictions of *Wolbachia* infection
137 status using illumina data have previously been shown to have 98.8% concordance with PCR based
138 predication of infection status³².

139 **Nucleotide diversity**

140 Levels of polymorphism for mtDNA and *Wolbachia* were estimated as π in 10 kb windows using
141 VCFtools (v0.1.14)⁶⁹. To investigate whether the frequency spectrum conformed to the standard neutral
142 model of molecular evolution I also calculated Tajima's *D* in 10 kb windows using VCFtools. To assess

143 the significance of deviations in Tajima's D and π 10,000 simulations were performed using msms
144 conditioned on the number of variable sites and with no recombination⁷⁰.

145 **Linkage disequilibrium**

146 Linkage between *Wolbachia* and mtDNA SNPs could potentially be a predictor of co-inheritance of
147 mtDNA and *Wolbachia*. Linkage was estimated using VCFtools (v0.1.14) using inter-chrom-geno-r2 to
148 estimate r^2 between each SNP in the two genomes⁶⁷.

149 **Estimation of mtDNA and *Wolbachia* copy number**

150 In insects, the phenotypic effect of *Wolbachia* will vary depending upon copy number in the host cells^{9,32}.
151 Given that there are two copies of autosomal DNA in a cell, I infer mtDNA and *Wolbachia* copy number
152 based on the ratio between mtDNA and autosomal DNA. This is intended to provide a relative estimate of
153 copy number rather than an absolute measure. Relative copy number estimated in this way obscures intra-
154 individual variation and variation between tissues, though the authors note that all flies used in
155 constructing the libraries were females of approximately the same age. *Wolbachia* contains several
156 regions which were excluded due to unusually high coverage across all samples (more than 3x the mean
157 coverage). Average coverage of autosomal DNA was calculated from randomly chosen and equivalently
158 sized nuclear regions for each mtDNA (Scf_3L:8000000..8014945) and *Wolbachia*
159 (Scf_2L:11000000..11445873). The average coverage of each nuclear region, respectively, was then used
160 to normalize estimates of copy number for each genotype. Previously the results of measuring *Wolbachia*
161 copy number in the same samples using both qPCR and estimates from illumina read depth had a
162 Pearson's correlation coefficient of .79, thus this is a robust approach to measuring *Wolbachia* titre³¹.

163 **Phylogenomic analysis**

164 To understand the relationship between *Wolbachia* infection and mtDNA I reconstructed the genealogical
165 history of each within the sample population. Multiple alignments were generated for both mtDNA and
166 *Wolbachia* by concatenating fasta consensus sequence files for each genotype. All indels and non-biallelic
167 SNPs were excluded from the dataset prior to generating the consensus fasta for each genotype. RAxML
168 version 8.10.2 was used to reconstruct phylogenies⁷¹. Maximum likelihood tree searches were conducted

169 using a general time reversible (GTR) model of nucleotide substitution with CAT rate heterogeneity and
170 all model parameters estimated by RAxM⁷². Trees were inferred using the rapid bootstrap algorithm and
171 simultaneous estimation of trees and bootstrapping, with automatic estimation of the necessary number of
172 bootstrap replicates.

173 **Association Analysis**

174 The association analysis focused on a relationship between nuclear polymorphisms and *Wolbachia* and
175 mtDNA copy number. To reduce the need for correction due to multiple testing and focus on regions that
176 may have been affected by selection due to the recent invasion of *Wolbachia* I used a subset of the
177 nuclear genome identified previously as exhibiting haplotype structure suggestive of recent selection^{62,73}.
178 These regions are unusually long haplotype blocks, thus many of the SNPs within each block are not
179 independent, reducing the need for correction due to multiple testing. Heterozygous bases were coded as
180 missing, and all loci with more than 10% missing data were excluded from the analysis, as well as SNPs
181 with a minor allele frequency of less than 2%, meaning they were present in the population in at least 3
182 copies. mtDNA and *Wolbachia* copy number were used for a multivariate analysis of association using
183 plink.multivariate⁷⁴. To investigate the possibility that *Wolbachia* copy number is affected by
184 polymorphisms in mtDNA, and vice versa, a single trait analysis was performed using plink v. 1.07⁷⁵.

185

186 RESULTS

187 **Sequencing Data**

188 The autosomal data included in this analysis was reported in⁶². There was very little variation in both
189 *Wolbachia* and mtDNA in this population. This included 78 SNPs and indels in the *Wolbachia* genome
190 and 90 in mtDNA. Reduced diversity has been reported previously in *D. simulans* mtDNA^{24,25}. The
191 authors note that previous work has established that there is no unusual relatedness in the nuclear genome
192 of this population⁶².

193 **Infection status**

194 In a previous study lines were scored as infected if they had a breadth of coverage greater than 90% and a
195 mean depth greater than one³². However, that dataset had a clearly bimodal distribution between infected
196 and uninfected lines, where uninfected lines had breadth of coverage less than 10% while infected lines
197 had a breadth of coverage of greater than 90%. As such that this demarcation was a natural interpretation
198 of the data³². In *D. simulans*, all lines had ~99% breadth of coverage aside from a single line with both a
199 lower overall depth of coverage and 80% breadth (Fig. 1). For this reason, all lines were scored as
200 infected. 100% infection is not unusually high for *D. simulans*.

201 **Nucleotide Diversity**

202 Estimates of π in *Wolbachia* ranged from 5.98×10^{-7} to 1×10^{-3} , with an average of 1.42×10^{-5} , within the
203 range of estimates from *Wolbachia* in *D. melanogaster* from another study ($7.9 \times 10^{-6} - 2.8 \times 10^{-5}$)³¹. The
204 mean of π in simulated populations of *Wolbachia* is 1.9×10^{-3} suggesting that variation is somewhat
205 reduced in wRi . π in mtDNA is 1×10^{-4} which again is similar to estimates from *D. melanogaster* ($4.34 \times$
206 $10^{-4} - 1.51 \times 10^{-3}$)³¹.

207 Overall Tajima's D was estimated to be -2.4 for *D. simulans* mtDNA (Fig 2). This is similar to
208 estimates in *D. melanogaster*³². Significance of this estimate was assessed using 10,000 simulations in
209 msms conditioned on the number of segregating sites and no recombination, and it is significant at p
210 $< .05$. Tajima's D in *Wolbachia* is not significantly different from expectations under neutrality based on
211 10,000 simulations. Thus, while a selective sweep or population bottleneck seems to have strongly
212 effected mtDNA in *D. simulans*, the same is not true of the *Wolbachia* population (Fig 2). This is very
213 different from *D. melanogaster* where *Wolbachia* and mtDNA had similar patterns of nucleotide
214 diversity³².

215 This is also much more negative than previously reported for mtDNA in *D. simulans*²⁵. It is very
216 different from the general patterns of Tajima's D in the nuclear genome, where average Tajima's D is 1
217 and the majority of the genome has a positive Tajima's D . Simulations in previous work suggest that the
218 pervasively positive values in the nuclear genome may be due to a population contraction, which again

219 indicates that the population dynamics affecting *D. simulans* nuclear and mtDNA genomes are very
220 different^{25,62}.

221 **Linkage disequilibrium**

222 There was no significant linkage disequilibrium between the genomes of *Wolbachia* and *D. simulans*
223 mtDNA. Average LD between *Wolbachia* and mtDNA SNPs was 2.06×10^{-3} . This may be because the
224 infection of *D. simulans* was too recent for variation to accumulate along particular lineages, and also
225 suggests that *D. simulans* was infected by a single invasion.

226

227 **Estimation of mtDNA and *Wolbachia* copy number**

228 There was considerable heterogeneity in both *Wolbachia* and mtDNA copy number (Fig. 1). Mean
229 (standard deviation) copy number of *Wolbachia* is 5.56 (2.45). This is similar to one estimate in *D.*
230 *melanogaster*, where mean copy number is 5.57 (3.95) though the standard deviation is lower in *D.*
231 *simulans*³². The reported mean was lower in other populations of *D. melanogaster*, though still within the
232 same range (2 - 4.5)³¹. Similarly mean mtDNA copy number is 33.85 (15.5) in *D. simulans* and 32.9
233 (44.5) in one estimate for *D. melanogaster*³². This is again not an absolute measure, but relative to nuclear
234 genomic coverage. The lower standard deviation could be due to more precise staging of the age of *D.*
235 *simulans*, less background variation effecting copy number (the *D. melanogaster* sample was from
236 multiple populations), or other unknown mechanisms. There was a positive relationship between mtDNA
237 and *Wolbachia* copy number (Fig. 1) ($p < 2.4 \times 10^{-7}$). While the functional reasons for or consequences of
238 this are unclear, because they are correlated they will be used in a multivariate analysis of association
239 rather than as separate analyses.

240 **Phylogenomic analysis**

241 To understand the relationship between *Wolbachia* infection status and mtDNA sequence variation I
242 reconstructed the phylogenetic history of the complete *Wolbachia* and mtDNA genome using the entire
243 set of 167 strains (Fig 3-4). What I found is consistent with the recent spread of *Wolbachia* in *D.*
244 *simulans*, as both phylogenies are essentially unresolved. This is not unexpected for mtDNA given

245 previous work in the species which found little within-haplotype variation among the three major mtDNA
246 haplotypes in *D. simulans*^{25,28}. Furthermore, of the 167 sequences 88 are identical to at least one other
247 sequence in the sample. While the *Wolbachia* phylogenetic tree gives the impression of having more
248 resolution than mtDNA, this is likely due to the larger genome, as the branches have similarly low
249 support. Of the 167 strains included in the tree 18 are identical to one or more *Wolbachia* genomes. Both
250 trees are essentially star phylogenies with the majority of bootstrap support values being less than 30.
251 Bootstrap support of greater than 70, for two branches in the mtDNA tree and five in the *Wolbachia* tree,
252 is shown (Fig 3-4). If uninfected individuals had been included in the dataset perhaps it would be possible
253 to test for congruence between the two phylogenies, however the essentially unresolved trees make it
254 clear that both *Wolbachia* and mtDNA swept the population recently.

255

256 **Association Analysis**

257 Association analysis was performed using `plink.multivariate` by regressing the line means for mtDNA and
258 *Wolbachia* copy number on each SNP contained within the previously identified in a scan for selection⁶².
259 This scan for selection focused on identifying haplotype blocks in LD. This considerably reduces the
260 number of SNPs tested for association, in addition to the fact that the SNPs are in haplotype blocks and
261 are therefore not independent tests^{62,73}. This reduces the need for correction due to multiple testing. I used
262 a *p*-value cut-off of $p < 9 \times 10^{-6}$ and identified 16 SNPs associated with *Wolbachia* and mtDNA copy
263 number. Of these 16 SNPs 13 are located in the same region on chromosome 2R (`Scf_2R: 13550916-`
264 `13569038`). Given the concentration of significant SNPs in a single region, this is also the region I will
265 focus on the most in the following discussion. The region containing 13 SNPs contains nine genes, four of
266 which are involved in CAAX-box protein processing, *ste24a-c* and a recent duplicate of *ste24c* CG30461.
267 CAAX-box protein processing is a part of a series of posttranslational protein modifications collectively
268 called protein prenylation which are required for fully functional proteins to be targeted to cell
269 membranes or organelles. It has been shown that pathogenic bacteria can exploit the host cell's
270 prenylation machinery, though it is unclear if this occurs in *Wolbachia*⁵⁹.

271 The other five genes are *AsnRs-m*, which is largely unannotated but is thought to a mitochondrial
272 aminoacyl-tRNA synthetase⁷⁶. *NIPP1Dm* is involved in axon guidance and negative regulation of protein
273 phosphorylation^{77,78}. *CG6805* is generally unannotated but is inferred to be involved in
274 dephosphorylation⁷⁶. *Cbp53E* regulates neural development⁷⁹. Lastly, *Ehbp1* is a developmental gene
275 implicated in regulation of the Notch pathway and membrane organization⁸⁰.

276 Of the other three SNPs identified in this association analysis two are located at Scf_2R:5814103
277 and Scf_2R:5811043, while the third is located at Scf_3L:2055556. Scf_2R:5811043 and
278 Scf_2R:5814103 are located in *Su(var)2-10* and *Phax*, respectively. These are neighboring genes, though
279 there is a third gene within 10 kb, *Mys45A*. *Su(var)2-10* is involved in development and chromosome
280 organization, but it has also been implicated in the regulation of the innate immune response and defense
281 against gram-negative bacteria⁸¹. *Su(var)2-10* is of particular interest given that *Wolbachia* are gram-
282 negative bacteria, however the potential role of *Su(var)2-10* in immune response is not clear. *Phax* is not
283 well annotated but is inferred to be involved in snRNA export from the nucleus⁷⁹. *Mys45A* is potentially
284 involved in actin cytoskeleton organization⁷⁹. In *D. melanogaster* *Wolbachia* uses host actin for maternal
285 transmission, though this has not been verified in *D. simulans*⁸². The last SNP, at Scf_3L:2055556, is in
286 *Connectin*, a cell adhesion protein also involved in axon guidance⁸³.

287 The identification of these SNPs in association with mtDNA and *Wolbachia* copy number does
288 not imply a functional relationship. Nonetheless, I chose to investigate whether any of these substitutions
289 had an effect on the coding sequence of any of genes in the region. Of the three SNPs found outside the
290 region containing the CAAX-box proteins all were either in introns or regulatory regions. Of the 13 SNPs
291 identified between Scf_2R: 13550916-13569038 eight are in introns or untranslated regions, including
292 one in the long intron of *Cb53E*, three in the introns or noncoding transcript of *CG6805*, and two in the
293 introns of *Ephb*. Of the remaining five SNPs four are in coding regions but silent, causing no change in
294 the amino acid sequence of the protein. This includes silent mutations in the exons of *ste24c* and two
295 silent mutations in the exons of *Ephb*. One SNP located in an exon of *ste24a*, at 13558515, is an amino
296 acid substitution from a Leucine to a Valine. This is not an uncommon amino acid substitution^{84,85}, though

297 it can be associated with phenotypes^{86,87}. Mutations in introns and untranslated regions could also be
298 having an effect on gene expression or processing, as could other linked SNPs in the region that were not
299 included in the analysis.

300

301 *Association between Wolbachia and mtDNA*

302 Association analysis was performed using plink by regressing the line means for mtDNA copy number
303 onto the *Wolbachia* genome and vice versa⁷⁵. There was no association between *Wolbachia* SNPs and
304 mtDNA copy number, but the opposite was not true. One SNP in the *D. simulans* mtDNA affected
305 *Wolbachia* copy number at $p < 3.18 \times 10^{-6}$. It is located in the *D. simulans* homolog of *D. melanogaster*
306 *srRNA* which has been implicated in pole cell formation⁸⁸. *Wolbachia* is incorporated into the pole cells,
307 the precursor to the germline, in order to be transmitted⁸⁹.

308

309 DISCUSSION

310 Using high through-put sequencing of a large panel of *D. simulans* I have reconstructed the complete
311 genome sequences of mtDNA and *Wolbachia*. I use these genome sequences to investigate the recent
312 history of *Wolbachia* and mtDNA in this population, as well as to estimate titre of both *Wolbachia* and
313 mtDNA. The history of *Wolbachia* in this population is reflected in the essentially star-like phylogeny of
314 both mtDNA and *Wolbachia*, indicating recent spread and co-inheritance. Lack of variation at mtDNA
315 and *Wolbachia* suggests a single spread of *w*Ri in this population as well as strict vertical transmission in
316 the maternal cytoplasm. Variation in *Wolbachia* is within the range expected under a neutral model,
317 however that was not the case for mtDNA which suggests either a selection sweep or a population
318 bottleneck. Previous studies found similar population genetic patterns at *Wolbachia* and mtDNA in *D.*
319 *melanogaster*, and thus could not distinguish whether selection on *Wolbachia* was driving similar patterns
320 in mtDNA or vice versa³¹. The much stronger pattern of negative Tajima's *D* in the mtDNA suggests that
321 in *D. simulans* selection is in fact mitochondrial. There was no linkage disequilibrium between *Wolbachia*

322 and mtDNA variants, however this is most likely due to fixation of a single mitochondrial haplotype
323 without considerable subsequent mutation.

324 Currently little is known about how *Wolbachia* interacts with its host^{37-39,82,90}. Understanding these
325 interactions, including regulation of *Wolbachia* titre, will be key to understanding the evolution of
326 *Wolbachia* and its hosts. By normalizing *Wolbachia* and mtDNA copy number using coverage of the
327 nuclear genome I am able to obtain estimates of its abundance. Much as in previous work, mtDNA copy
328 number was higher than *Wolbachia* copy number, though both varied across strains³². As all of my data
329 was produced from adult females, at the same time, using the same techniques, there is no danger that this
330 is due to differences in methodology among samples³². Estimates of copy number were very similar to
331 previous work in *D. melanogaster*, performed with qPCR, and there has been shown to be a high
332 correlation between qPCR and illumina estimates of copy number^{31,32}. These are not absolute measures,
333 rather they are relative to one another and to nuclear copy number, and they provide robust estimates of
334 *Wolbachia* titre within the population. As the *Wolbachia* phylogenetic tree is essentially unresolved in
335 this population but there is considerable variation in *Wolbachia* titre, it is clear that some host factors
336 must be affecting variation in *Wolbachia* titre.

337 The history of mtDNA and the nuclear genome is quite divergent in this population. The nuclear
338 genome has an average Tajima's *D* of 1 and 5 polymorphisms for every 100 bp⁶². Simulations suggest
339 that this is due to a combination of population contraction and selection, most likely from standing
340 variation, though many types of sweeps can produce similar signatures⁶². In contrast the mtDNA genome
341 contains an abundance of low frequency variation, and in fact many of the mtDNA genomes sampled in
342 this population are identical. This is consistent with the recent spread, single origin, and maternal
343 transmission, of *wRi* in *D. simulans*. This is consistent with previous work which found low levels of
344 mtDNA variation in *D. simulans* within a haplotype^{24,91}. This is also consistent with work on *Wolbachia*
345 which documented the spread of *wRi* in *D. simulans* in the 1980's^{4,5,8,18-20}.

346 While it has been proposed elsewhere, the author is not aware of another association analysis of
347 *Wolbachia* and mtDNA copy number³². *Wolbachia* copy number is known to be affected by host

348 background, but the genes or mechanisms involved are not known^{54,55,57}. The fact that four of the nine
349 genes found in the primary region detected in the association analysis are involved in CAAX-box protein
350 processing is of particular interest, given the history of this type of gene and intracellular pathogens.
351 CAAX-box protein processing is a part of a series of posttranslational protein modifications collectively
352 called protein prenylation which are required for fully functional proteins to be targeted to cell
353 membranes or organelles. Prenylated proteins include Ras, Rac, and Rho. However, it has been shown
354 that pathogenic bacteria can exploit the host cell's prenylation machinery⁵⁹. For example, *Salmonella-*
355 *induced filament A* is a protein from *Salmonella typhimurium*, a gram-negative facultative intracellular
356 bacterium. *Salmonella-induced filament A* has a CAAX motif required for prenylation to occur, it was
357 shown to be processed by host prenylation machinery, and it is necessary for survival of the
358 bacterium^{60,92,93}. *Legionella pneumophila* Ankyrin B protein exploits the host prenylation machinery in
359 order to anchor Ankyrin B protein to the membrane of the pathogenic vacuole⁶¹. Proliferation of
360 *Legionella pneumophila* requires Ankyrin B, as does the manifestation of Legionnaires disease. Ankyrin
361 repeat domains are most commonly found in eukaryotes and viruses, though they are rarely found in
362 bacteria and Archaea⁹⁴. In bacteria they are found in a few obligate or facultative intracellular
363 Proteobacteria⁵⁹. *Wolbachia* has an unusually high number of Ankyrin repeat domains with rapid
364 evolution⁹⁴. Ankyrin proteins play a major role in host-pathogen interactions and the evolution of
365 infections^{95,96}. There is no way to know from the current analysis if the Ankyrin repeat genes are
366 exploiting the host prenylation system but it is an intriguing area for future investigation. The results of
367 this association analysis suggest that some interaction between the pathogen and its host is targeting the
368 protein prenylation machinery.

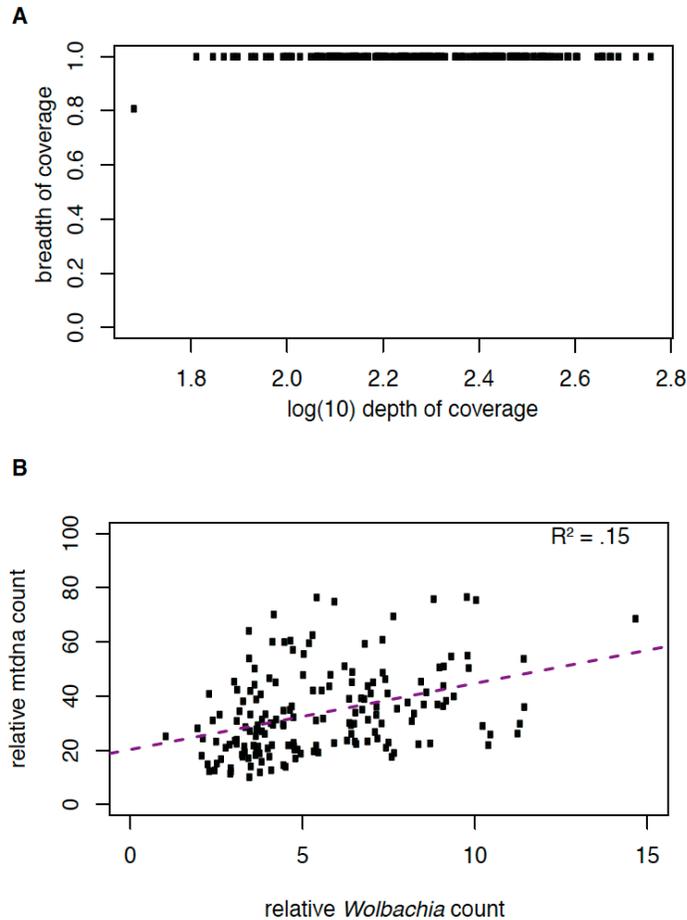
369 There was also an association between a polymorphism in *srRNA*, which has been implicated in
370 pole cell formation⁸⁸, and *Wolbachia* copy number. Concentration of *Wolbachia* in the posterior of the
371 embryo, where pole cells are forming, is correlated with degree of cytoplasmic incompatibility⁹⁷. *D.*
372 *simulans* has been shown to have nearly complete cytoplasmic incompatibility, though it is possible there
373 are mutations sorting at low frequency that affect this or that mitigate negative phenotypic consequences

374 of high *Wolbachia* titre. It has also been demonstrated that *gurken* is important for *Wolbachia* titre in the
375 germline in *D. melanogaster*, and it is involved in pole cell formation beginning at an earlier stage than
376 *srRNA* suggesting there could be an interaction between the two factors^{88,90}. *D. simulans* wRi has a
377 different distribution in the cytoplasm from other strains of *Wolbachia*, as it tends to evenly distribute
378 throughout the embryo while other strains are either concentrated at the posterior, or at the anterior of the
379 embryo away from the pole cells⁹⁷. Future work in related species may show that these different
380 distributions also mitigate different interactions between host and symbiont, including being effected by
381 different genes and processes within the host.

382

383

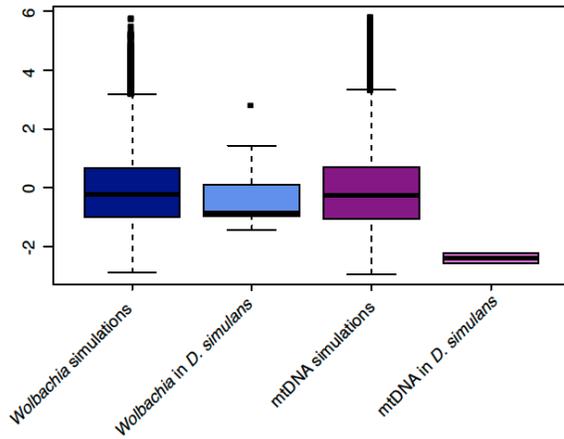
384



385

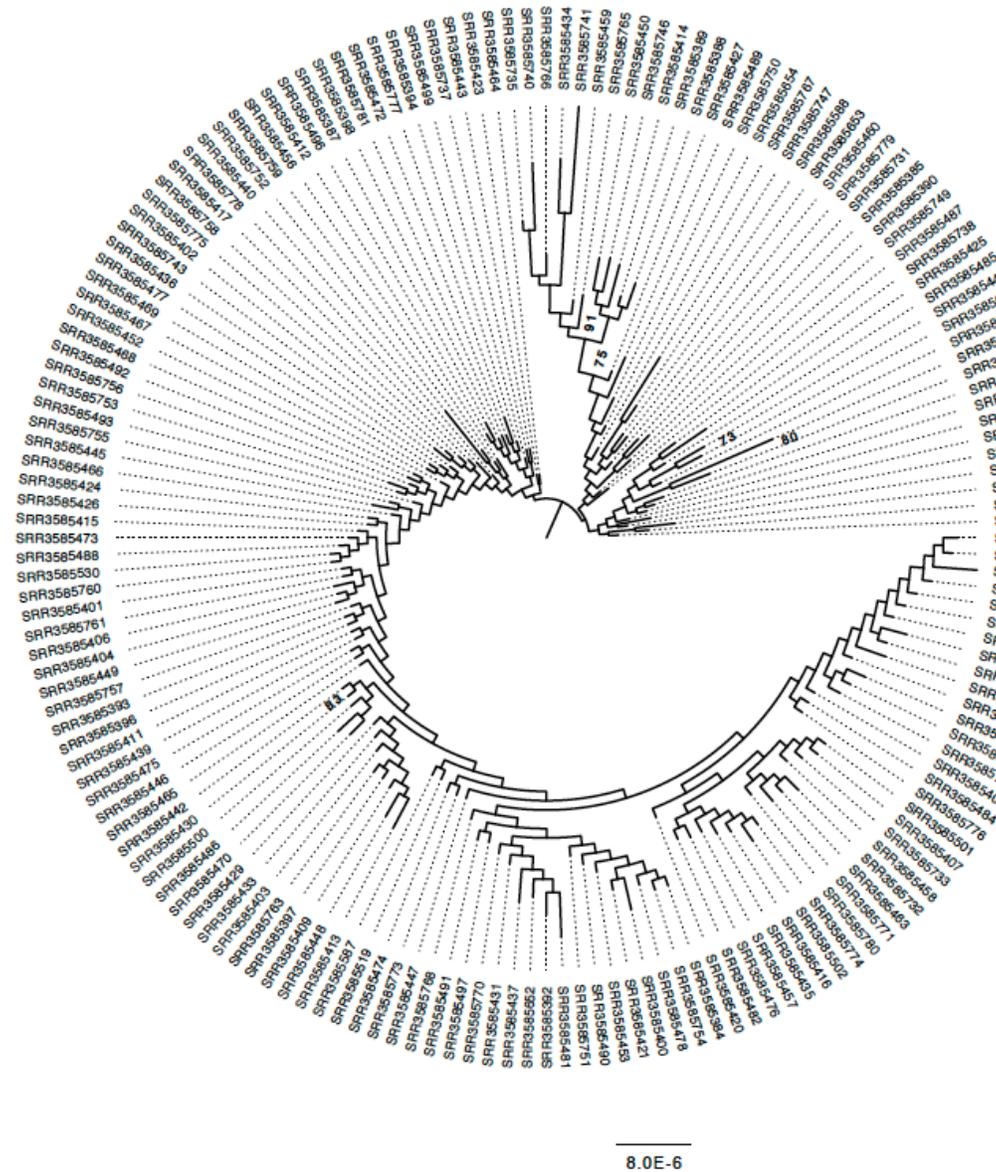
Figure 1. *Wolbachia* infection status and relationship to mtDNA copy number **A.** Relationship between depth and breadth of sequencing coverage for *Wolbachia* assemblies in the *D. simulans* panel. Depth of coverage is shown in \log_{10} unites and is calculated as the number of reads present at each nucleotide in the reference averaged over every site. Breadth of coverage is the proportion of covered nucleotides in the consensus sequence relative to the reference. **B.** Relationship between relative mtDNA copy number and *Wolbachia* copy number. Both were normalized relative to nuclear coverage. Although separate regions were used to normalize mtDNA and *Wolbachia*, as they are different sizes, average values were very similar within genotypes. The relationship between mtDNA and *Wolbachia* copy number is positive ($p < 2.4 \times 10^{-7}$).

386



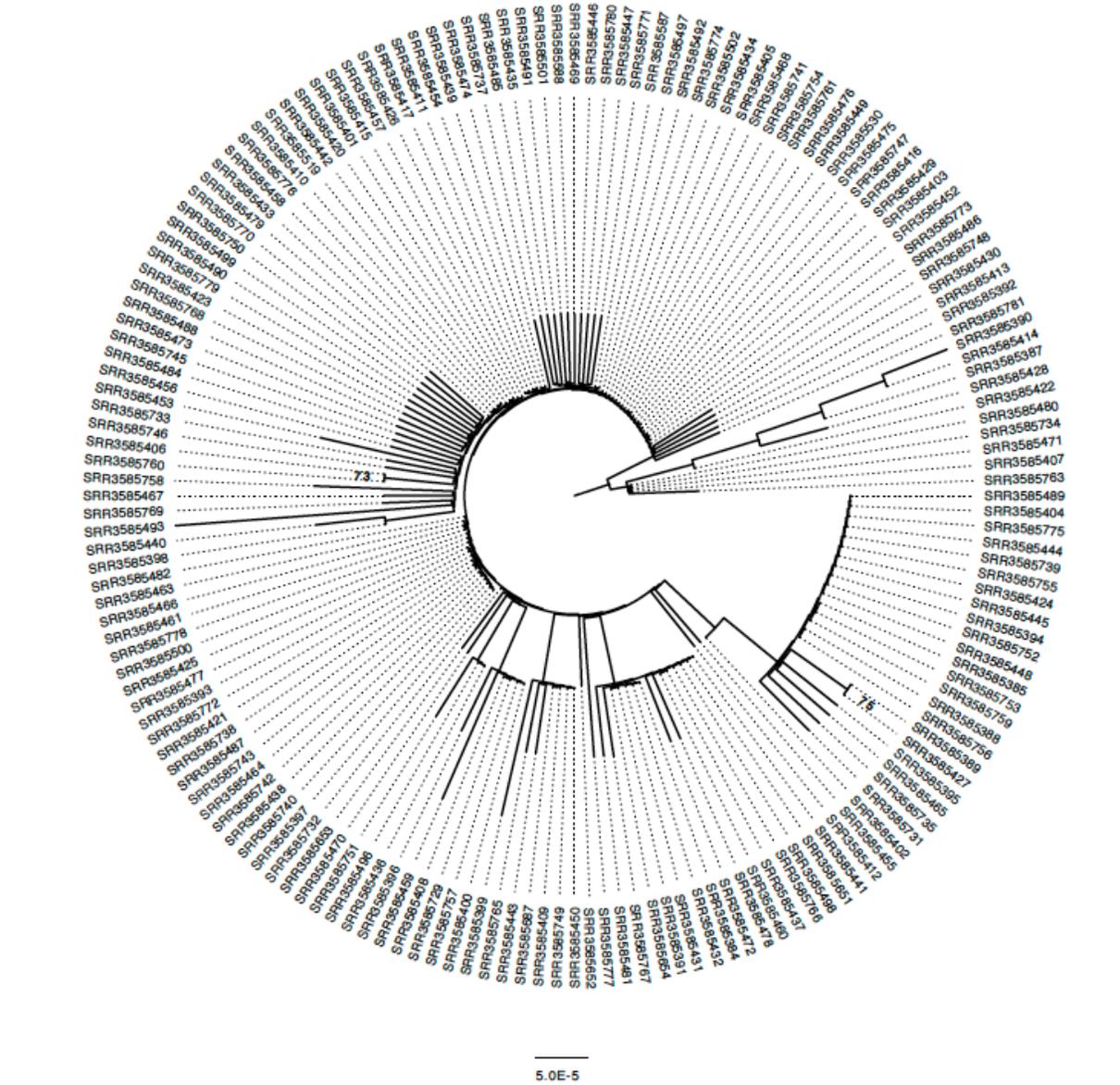
387

388 **Figure 2. *Wolbachia* and mtDNA Tajima's *D*.** 10,000 simulations were performed for *Wolbachia* and
389 *D. simulans* each conditioned upon the number of polymorphisms. The actual values in *D. simulans*
390 mtDNA are outside the 95% confidence interval of the simulations, while *Wolbachia* is not. There is
391 considerable variation in Tajima's *D* across the *Wolbachia* genome while mtDNA is much smaller and
392 invariant in its values of Tajima's *D*.



393

Figure 3: Maximum likelihood genealogy of the *D. simulans* *Wolbachia* pathogen. All strains were infected with *Wolbachia* and are included in this genealogy. The underlying data consist of an ungapped multiple alignment of 168 sequences of the entire *Wolbachia* genome. The unrooted tree was midpoint rooted for visualization and branches with > 70% RAxML bootstrap support values are shown in bold. Scale bars for branch lengths are in term of mutations per site. The majority of branches are essentially unsupported by bootstrapping.



394

Figure 4: Maximum likelihood genealogy of the *D. simulans* mtDNA genome. The underlying data consist of an ungapped multiple alignment of 168 sequences of the entire mtDNA genome. The unrooted tree was midpoint rooted for visualization and branches with > 70% RAxML bootstrap support values are shown in bold. Scale bars for branch lengths are in term of mutations per site. The tree is largely unresolved, suggesting recent spread of this mtDNA haplotype through the population.

395

396

397

398

399

1. Zug, R. & Hammerstein, P. Still a host of hosts for *Wolbachia*: analysis of recent data suggests that 40% of terrestrial arthropod species are infected. *PLoS ONE* 7, e38544 (2012).
2. Hilgenboecker, K., Hammerstein, P., Schlattmann, P., Telschow, A. & Werren, J. H. How

- 400 many species are infected with *Wolbachia*?--A statistical analysis of current data. *FEMS*
401 *Microbiol. Lett.* **281**, 215–220 (2008).
- 402 3. Jeyaprakash, A. & Hoy, M. A. Long PCR improves *Wolbachia* DNA amplification: wsp
403 sequences found in 76% of sixty-three arthropod species. *Insect Mol. Biol.* **9**, 393–405
404 (2000).
- 405 4. Turelli, M. & Hoffmann, A. A. Cytoplasmic incompatibility in *Drosophila simulans*:
406 dynamics and parameter estimates from natural populations. *Genetics* **140**, 1319–1338
407 (1995).
- 408 5. Turelli, M., Hoffmann, A. A. & McKechnie, S. W. Dynamics of cytoplasmic
409 incompatibility and mtDNA variation in natural *Drosophila simulans* populations. *Genetics*
410 **132**, 713–723 (1992).
- 411 6. Charlat, S., Nirgianaki, A., Bourtzis, K. & Merçot, H. Evolution of *Wolbachia*-induced
412 cytoplasmic incompatibility in *Drosophila simulans* and *D. sechellia*. *Evolution* **56**, 1735–
413 1742 (2002).
- 414 7. Rousset, F., Vautrin, D. & Solignac, M. Molecular identification of *Wolbachia*, the agent of
415 cytoplasmic incompatibility in *Drosophila simulans*, and variability in relation with host
416 mitochondrial types. *Proc. R. Soc. B* **247**, 163–168 (1992).
- 417 8. Hoffmann, A. A., Turelli, M. & Harshman, L. G. Factors affecting the distribution of
418 cytoplasmic incompatibility in *Drosophila simulans*. *Genetics* **126**, 933–948 (1990).
- 419 9. Berticat, C., Rousset, F., Raymond, M., Berthomieu, A. & Weill, M. High *Wolbachia*
420 density in insecticide-resistant mosquitoes. *Proc. R. Soc. B* **269**, 1413–1416 (2002).
- 421 10. Wong, Z. S., Hedges, L. M., Brownlie, J. C. & Johnson, K. N. *Wolbachia*-mediated
422 antibacterial protection and immune gene regulation in *Drosophila*. *PLoS ONE* **6**, e25430–9
423 (2011).
- 424 11. Osborne, S. E., Leong, Y. S., O'Neill, S. L. & Johnson, K. N. Variation in antiviral
425 protection mediated by different *Wolbachia* strains in *Drosophila simulans*. *PLoS Pathog* **5**,
426 e1000656–9 (2009).
- 427 12. Fytrou, A., Schofield, P. G., Kraaijeveld, A. R. & Hubbard, S. F. *Wolbachia* infection
428 suppresses both host defence and parasitoid counter-defence. *Proc. R. Soc. B* **273**, 791–796
429 (2006).
- 430 13. Brownlie, J. C. *et al.* Evidence for metabolic provisioning by a common invertebrate
431 endosymbiont, *Wolbachia pipientis*, during periods of nutritional stress. *PLoS Pathog* **5**,
432 e1000368 (2009).
- 433 14. Teixeira, L., Ferreira, Á. & Ashburner, M. The bacterial symbiont *Wolbachia* induces
434 resistance to RNA viral infections in *Drosophila melanogaster*. *PLoS Biol.* **6**, e1000002
435 (2008).
- 436 15. Hedges, L. M., Brownlie, J. C. & O'Neill, S. L. *Wolbachia* and virus protection in insects.
437 **322**, 702 *Science* (2008).
- 438 16. Jiggins, F. M., Hurst, G. & Majerus, M. Sex-ratio-distorting *Wolbachia* causes sex-role
439 reversal in its butterfly host B. *Proc. R. Soc. B* **267**, 69–73.
- 440 17. Koukou, K. *et al.* Influence of antibiotic treatment and *Wolbachia* curing on sexual isolation
441 among *Drosophila melanogaster* cage populations. *Evolution* **60**, 87–11 (2006).
- 442 18. Turelli, M. & Hoffmann, A. A. Rapid spread of an inherited incompatibility factor in
443 California *Drosophila*. *Nature* **353**, 440–442 (1991).
- 444 19. Hoffmann, A. A. & Turelli, M. Unidirectional incompatibility in *Drosophila simulans*:
445 inheritance, geographic variation and fitness effects. *Genetics* **119**, 435–444 (1988).
- 446 20. Hoffmann, A. A., Turelli, M. & Simmons, G. M. Unidirectional incompatibility between
447 populations of *Drosophila simulans*. *Evolution* **40**, 692–701 (1986).
- 448 21. Klasson, L. *et al.* The mosaic genome structure of the *Wolbachia* wRi strain infecting
449 *Drosophila simulans*. *Proc. Nat. Acad. Sci. USA* **106**, 5725–5730 (2009).
- 450 22. Choi, J. Y. & Aquadro, C. F. The coevolutionary period of *Wolbachia pipientis* infecting

- 451 *Drosophila ananassae* and its impact on the evolution of the host germline stem cell
452 regulating genes. *Mol. Biol. Evol.* **31**, 2457–2471 (2014).
- 453 23. Ballard, J. W. O. Comparative genomics of mitochondrial DNA in *Drosophila simulans*. *J*
454 *Mol Evol* **51**, 64–75 (2000).
- 455 24. Ballard, J. W., Hatzidakis, J., Karr, T. L. & Kreitman, M. Reduced variation in *Drosophila*
456 *simulans* mitochondrial DNA. *Genetics* **144**, 1519–1528 (1996).
- 457 25. Ballard, J. W. O. Comparative genomics of mitochondrial DNA in *Drosophila simulans*. *J*
458 *Mol Evol* **51**, 64–75 (2000).
- 459 26. Solignac, M., Vautrin, D. & Rousset, F. Widespread occurrence of the proteobacteria
460 *Wolbachia* and partial cytoplasmic incompatibility in *Drosophila melanogaster*. *C. R. Acad.*
461 *Sci.* **317**, 461–479 (1994).
- 462 27. Solignac, M. & Monnerot, M. Race formation, speciation, and introgression within
463 *Drosophila simulans*, *D. mauritiana*, and *D. sechellia* inferred from mitochondrial DNA
464 analysis. *Evolution* **40**, 531–539 (1986).
- 465 28. Baba-Aïssa, F., Solignac, M., Dennebouy, N. & David, J. R. Mitochondrial DNA variability
466 in *Drosophila simulans*: quasi absence of polymorphism within each of the three
467 cytoplasmic races. *Heredity* **61**, 419–426 (1988).
- 468 29. Montchamp-Moreau, C., Ferveur, J. F. & Jacques, M. Geographic distribution and
469 inheritance of three cytoplasmic incompatibility types in *Drosophila simulans*. *Genetics*
470 **129**, 399–407 (1991).
- 471 30. James, A. C. & Ballard, J. W. Expression of cytoplasmic incompatibility in *Drosophila*
472 *simulans* and its impact on infection frequencies and distribution of *Wolbachia pipientis*.
473 *Evolution* **54**, 1661–1672 (2000).
- 474 31. Early, A. M. & Clark, A. G. Monophyly of *Wolbachia pipientis* genomes within *Drosophila*
475 *melanogaster*: geographic structuring, titre variation and host effects across five
476 populations. *Mol. Ecol.* **22**, 5765–5778 (2013).
- 477 32. Richardson, M. F. *et al.* Population genomics of the *Wolbachia* endosymbiont in *Drosophila*
478 *melanogaster*. *PLoS Genet.* **8**, e1003129 (2012).
- 479 33. Nunes, M. D. S., Nolte, V. & Schlatterer, C. Nonrandom *Wolbachia* infection status of
480 *Drosophila melanogaster* strains with different mtDNA haplotypes. *Mol. Biol. Evol.* **25**,
481 2493–2498 (2008).
- 482 34. Beckmann, J. F., Ronau, J. A. & Hochstrasser, M. A *Wolbachia* deubiquitylating enzyme
483 induces cytoplasmic incompatibility. *Nat Microbiol* **2**, 17007 (2017).
- 484 35. Civetta, A. Rapid Evolution and gene-specific patterns of selection for three genes of
485 spermatogenesis in *Drosophila*. *Mol. Biol. Evol.* **23**, 655–662 (2005).
- 486 36. Bauer DuMont, V. L., Flores, H. A., Wright, M. H. & Aquadro, C. F. Recurrent positive
487 selection at *bgen*, a key determinant of germ line differentiation, does not appear to be
488 driven by simple coevolution with its partner protein *bam*. *Mol. Biol. Evol.* **24**, 182–191
489 (2006).
- 490 37. Flores, H. A., Bunnell, J. E., Aquadro, C. F. & Barbash, D. A. The *Drosophila bag of*
491 *marbles* gene interacts genetically with *Wolbachia* and shows female-specific effects of
492 divergence. *PLoS Genet.* **11**, e1005453 (2015).
- 493 38. Serbus, L. R. & Sullivan, W. A Cellular basis for *Wolbachia* recruitment to the host
494 germline. *PLoS Pathog* **3**, e190 (2007).
- 495 39. Serbus, L. R., Casper-Lindley, C., Landmann, F. & Sullivan, W. The genetics and cell
496 biology of *Wolbachia*-host interactions. *Annu. Rev. Genet.* **42**, 683–707 (2008).
- 497 40. McGraw, E. A., Merritt, D. J., Droller, J. N. & O'Neill, S. L. *Wolbachia* density and
498 virulence attenuation after transfer into a novel host. *Proc. Nat. Acad. Sci. USA* **99**, 2918–
499 2923 (2002).
- 500 41. Bordenstein, S. R., Marshall, M. L., Fry, A. J., Kim, U. & Wernegreen, J. J. The tripartite
501 associations between Bacteriophage, *Wolbachia*, and Arthropods. *PLoS Pathog* **2**, e43–10

- 502 (2006).
- 503 42. Martinez, J. *et al.* Should symbionts be nice or selfish? Antiviral effects of *Wolbachia* are
504 costly but reproductive parasitism is not. *PLoS Pathog* **11**, e1005021 (2015).
- 505 43. Clark, M. E., Veneti, Z., Bourtzis, K. & Karr, T. L. *Wolbachia* distribution and cytoplasmic
506 incompatibility during sperm development: the cyst as the basic cellular unit of CI
507 expression. *Mech. Dev.* **120**, 185–198 (2003).
- 508 44. Chrostek, E., Marialva, M. S. P., Yamada, R., O'Neill, S. L. & Teixeira, L. High anti-viral
509 protection without immune upregulation after interspecies *Wolbachia* transfer. *PLoS ONE* **9**,
510 e99025 (2014).
- 511 45. Chrostek, E. & Teixeira, L. Mutualism breakdown by amplification of *Wolbachia* genes.
512 *PLoS Biol.* **13**, e1002065 (2015).
- 513 46. Chrostek, E. *et al.* *Wolbachia* variants induce differential protection to viruses in *Drosophila*
514 *melanogaster*: A phenotypic and phylogenomic analysis. *PLoS Genet.* **9**, e1003896 (2013).
- 515 47. Hoffmann, A. A. *et al.* Successful establishment of *Wolbachia* in *Aedes* populations to
516 suppress dengue transmission. *Nature* **476**, 454–457 (2011).
- 517 48. McMeniman, C. J. *et al.* Host adaptation of a *Wolbachia* strain after long-term serial
518 passage in mosquito cell lines. *App. Environ. Microbiol.* **74**, 6963–6969 (2008).
- 519 49. Osborne, S. E., Iturbe-Ormaetxe, I., Brownlie, J. C., O'Neill, S. L. & Johnson, K. N.
520 Antiviral protection and the importance of *Wolbachia* density and tissue tropism in
521 *Drosophila simulans*. *App. Environ. Microbiol.* **78**, 6922–6929 (2012).
- 522 50. Kondo, N., Shimada, M. & Fukatsu, T. Infection density of *Wolbachia* endosymbiont
523 affected by co-infection and host genotype. *Biol. Lett.* **1**, 488–491 (2005).
- 524 51. Reynolds, K. T., Thomson, L. J. & Hoffmann, A. A. The effects of host age, host nuclear
525 background and temperature on phenotypic effects of the virulent *Wolbachia* strain popcorn
526 in *Drosophila melanogaster*. *Genetics* **164**, 1027–1034 (2003).
- 527 52. Olsen, K., Reynolds, K. T. & Hoffmann, A. A. A field cage test of the effects of the
528 endosymbiont *Wolbachia* on *Drosophila melanogaster*. *Heredity* (2001).
- 529 53. Boyle, L., O'Neill, S. L., Robertson, H. M. & Karr, T. L. Interspecific and intraspecific
530 horizontal transfer of *Wolbachia* in *Drosophila*. *Science* **260**, 1796–1799 (1993).
- 531 54. Poinot, D., Bourtzis, K., Markakis, G. & Savakis, C. *Wolbachia* transfer from *Drosophila*
532 *melanogaster* into *D. simulans*: host effect and cytoplasmic incompatibility relationships.
533 *Genetics* (1998).
- 534 55. McGraw, E. A., Merritt, D. J., Droller, J. N. & O'Neill, S. L. *Wolbachia*-mediated sperm
535 modification is dependent on the host genotype in *Drosophila*. *Proc. R. Soc. B* **268**, 2565–
536 2570 (2001).
- 537 56. Poinot, D., Montchamp-Moreau, C. & Mercot, H. *Wolbachia* segregation rate in
538 *Drosophila simulans* naturally bi-infected cytoplasmic lineages. *Heredity* **85 (Pt 2)**, 191–
539 198 (2000).
- 540 57. Boyle, L., O'Neill, S. L., Robertson, H. M. & Karr, T. L. Interspecific and intraspecific
541 horizontal transfer of *Wolbachia* in *Drosophila*. *Science* **260**, 1796–1799 (1993).
- 542 58. Perrot-Minnot & Werren. *Wolbachia* infection and incompatibility dynamics in
543 experimental selection lines. *J. Evol. Biol.* **12**, 272–282 (1999).
- 544 59. Amaya, M., Baranova, A. & van Hoek, M. L. Protein prenylation: A new mode of host-
545 pathogen interaction. *Biochemical and Biophysical Research Communications* **416**, 1–6
546 (2011).
- 547 60. Reinicke, A. T. *et al.* A *Salmonella typhimurium* effector protein SifA is modified by host
548 cell prenylation and s-acylation machinery. *J. Biol. Chem.* **280**, 14620–14627 (2005).
- 549 61. Price, C. T. D., Al-Quadani, T., Santic, M., Jones, S. C. & Abu Kwaiq, Y. Exploitation of
550 conserved eukaryotic host cell farnesylation machinery by an F-box effector of *Legionella*
551 *pneumophila*. *J Exp Med* **207**, 1713–1726 (2010).
- 552 62. Signor, S. A., New, F. & Nuzhdin, S. *in review*. An abundance of high frequency variance

- 553 uncovered in a large panel of *Drosophila simulans*.
554 63. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM
555 *arXiv*. 1–3 (2015).
556 64. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–
557 2079 (2009).
558 65. Meiklejohn, C. D. *et al.* An Incompatibility between a mitochondrial tRNA and its nuclear-
559 encoded tRNA synthetase compromises development and fitness in *Drosophila*. *PLoS*
560 *Genet.* **9**, e1003238 (2013).
561 66. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing
562 next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
563 67. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158
564 (2011).
565 68. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic
566 features. *Bioinformatics* **26**, 841–842 (2010).
567 69. Tajima, F. Statistical method for testing the neutral mutation hypothesis by DNA
568 polymorphism. *Genetics* **123**, 585–595 (1989).
569 70. Ewing, G. & Hermisson, J. MSMS: a coalescent simulation program including
570 recombination, demographic structure and selection at a single locus. *Bioinformatics* **26**,
571 2064–2065 (2010).
572 71. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
573 phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
574 72. Izquierdo-Carrasco, F., Smith, S. A. & Stamatakis, A. Algorithms, data structures, and
575 numerics for likelihood-based phylogenetic inference of huge trees. *BMC Bioinformatics*
576 **12**, 470 (2011).
577 73. Garud, N. R., Messer, P. W., Buzbas, E. O. & Petrov, D. A. Recent selective sweeps in
578 North American *Drosophila melanogaster* show signatures of soft sweeps. *PLoS Genet.* **11**,
579 e1005004 (2015).
580 74. Ferreira, M. A. R. & Purcell, S. M. A multivariate test of association. *Bioinformatics* **25**,
581 132–133 (2008).
582 75. Purcell, S. *et al.* PLINK: A tool set for whole-genome association and population-based
583 linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
584 76. FlyBase Curators, Swiss-Prot Project Members InterPro Project Members. Gene Ontology
585 annotation in FlyBase through association of InterPro records with GO terms. (2004).
586 77. Babu, K., Bahri, S., Alpey, L. & Chia, W. Bifocal and PP1 interaction regulates targeting
587 of the R-cell growth cone in *Drosophila*. *Dev. Biol.* **288**, 372–386 (2005).
588 78. Bennett, D., Szoor, B., Gross, S., Vereshchagina, N. & Alpey, L. Ectopic expression of
589 inhibitors of Protein phosphatase type 1 (PP1) can be used to analyze roles of PP1 in
590 *Drosophila* development. *Genetics* **164**, 235–245 (2003).
591 79. Gaudet, P., Livstone, M. & Thomas, P. Gene Ontology annotation inferences using
592 phylogenetic trees. GO Reference Genome Project. (2010).
593 80. Giagtzoglou, N. *et al.* *dEHBPI* controls exocytosis and recycling of *Delta* during
594 asymmetric divisions. *J. Cell Biol.* **196**, 65–83 (2012).
595 81. Cronin, S. J. F. *et al.* Genome-wide RNAi screen identifies genes involved in intestinal
596 pathogenic bacterial infection. *Science* **325**, 340–343 (2009).
597 82. Newton, I. L. G., Savytskyy, O. & Sheehan, K. B. *Wolbachia* utilize host actin for efficient
598 maternal transmission in *Drosophila melanogaster*. *PLoS Pathog* **11**, e1004798 (2015).
599 83. Chiba, A. Early development of the *Drosophila* neuromuscular junction: a model for
600 studying neuronal networks in development. *Int. Rev. Neurobiol.* (1999).
601 84. Yampolsky, L. Y. & Stoltzfus, A. The exchangeability of amino acids in proteins. *Genetics*
602 **170**, 1459–1472 (2005).
603 85. Creixell, P., Schoof, E. M., Tan, C. S. H. & Lindig, R. Mutational properties of amino acid

- 604 residues: implications for evolvability of phosphorylatable residues. *Philos. Trans. R. Soc. B*
605 **367**, 2584–2593 (2012).
- 606 86. Leone, A. *et al.* Evidence for nm23 RNA overexpression, DNA amplification and mutation
607 in aggressive childhood neuroblastomas. *Oncogene* **8**, 855–865 (1993).
- 608 87. Ishiko, A. *et al.* A novel leucine to valine mutation in residue 7 of the helix initiation motif
609 of Keratin10 leads to bullous congenital ichthyosiform erythroderma. *J. Inv. Dermatol.* **116**,
610 991–992 (2001).
- 611 88. Mahowald, A. P. Assembly of the *Drosophila* germ plasm. *Int. Rev. Cytol.* **203**, 187–213
612 (2001).
- 613 89. Kose, H. & Karr, T. L. Organization of *Wolbachia pipientis* in the *Drosophila* fertilized egg
614 and embryo revealed by an anti-*Wolbachia* monoclonal antibody. *Mech. Dev.* **51**, 275–288
615 (1995).
- 616 90. Serbus, L. R. *et al.* A feedback loop between *Wolbachia* and the *Drosophila* *gurken* mRNP
617 complex influences *Wolbachia* titer. *J. Cell Sci.* **124**, 4299–4308 (2012).
- 618 91. Ballard, J. W. O. Comparative genomics of mitochondrial DNA in *Drosophila simulans*. *J.*
619 *Mol. Evol.* **51**, 64–75 (2000).
- 620 92. Vinh, D. B. N., Ko, D. C., Rachubinski, R. A., Aitchison, J. D. & Miller, S. I. Expression of
621 the *Salmonella* *spp.* virulence factor SifA in yeast alters Rho1 activity on peroxisomes. *Mol.*
622 *Biol. Cell* **21**, 3567–3577 (2010).
- 623 93. Brumell, J. H., Goosney, D. L. & Finlay, B. B. SifA, a type III secreted effector of
624 *Salmonella typhimurium*, directs *Salmonella*-induced filament (Sif) formation along
625 microtubules. *Traffic* **3**, 407–415 (2002).
- 626 94. Siozios, S. *et al.* The diversity and evolution of *Wolbachia* ankyrin repeat domain genes.
627 *PLoS ONE* **8**, e55390 (2013).
- 628 95. Bork, P. Hundreds of ankyrin-like repeats in functionally diverse proteins: mobile modules
629 that cross phyla horizontally? *Proteins* **17**, 363–374 (1993).
- 630 96. Habyarimana, F. *et al.* Role for the Ankyrin eukaryotic-like genes of *Legionella*
631 *pneumophila* in parasitism of protozoan hosts and human macrophages. *Environ. Microbiol.*
632 **10**, 1460–1474 (2008).
- 633 97. Veneti, Z., Clark, M. E., Karr, T. L., Savakis, C. & Bourtzis, K. Heads or Tails: Host-
634 Parasite Interactions in the *Drosophila-Wolbachia* System. *App. Env. Microbiol.* **70**, 5366–
635 5372 (2004).

636

637 **Acknowledgements**

638 I would like to thank J. Butler for helpful commentary on this manuscript, and S. Nuzhdin for helpful
639 discussion and direction.

640

641 **Author Contributions**

642 S. S. performed all of the work described herein.

643

644 **Competing financial interest**

645 The author has no competing financial interests.